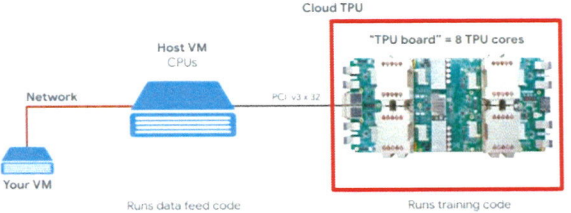
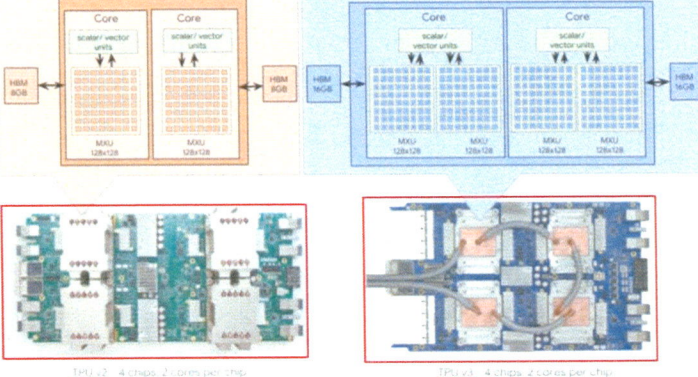
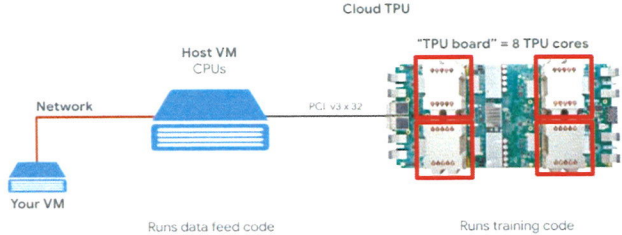
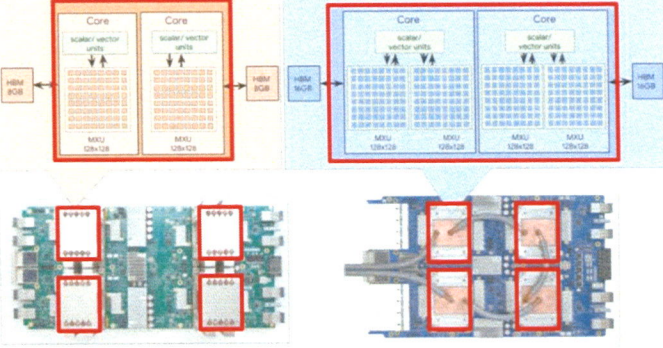


# EXHIBIT I

# **Exhibit B**

**U.S. Pat. No. 8,407,273**  
**Claim 53**

'273 PATENT	INFRINGEMENT EVIDENCE
<p>53. A <b>device</b>:  comprising at least one first low precision high-dynamic range (LPHDR) execution unit adapted to execute a first operation on a first input signal representing a first numerical value to produce a first output signal representing a second numerical value,  wherein the dynamic range of the possible valid inputs to the first operation is at least as wide as from 1/1,000,000 through 1,000,000 and for at least <math>X=5\%</math> of the possible valid inputs to the first operation, the statistical mean, over repeated execution of the first operation on each specific input from the at least <math>X\%</math> of the possible valid inputs to the first operation, of the numerical values represented by the first output signal of the LPHDR unit executing the first operation on that input differs by at least <math>Y=0.05\%</math> from the result of an exact mathematical calculation of the first operation on the numerical values of that same input;  wherein the number of LPHDR execution units in the device exceeds by at least one hundred the non-negative integer number of execution units in the device adapted to execute at least the operation of multiplication on floating point numbers that are at least 32 bits wide.</p>	<p>As demonstrated below, the Accused Products include multiple components that, separately and independently, meet all the requirements of the claimed "device." For example, a "TPU Board" satisfies these requirements:</p> <div data-bbox="961 272 1696 695"> <p>When you request one "Cloud TPU v2" on Google Cloud Platform, you get a virtual machine (VM) which has a PCI-attached TPU board. The TPU board has four dual-core TPU chips. Each TPU core features a VPU (Vector Processing Unit) and a 128x128 MXU (Matrix multiply Unit). This "Cloud TPU" is then usually connected through the network to the VM that requested it. So the full picture looks like this:</p>  <p>Illustration: your VM with a network-attached "Cloud TPU" accelerator. "The Cloud TPU" itself is made of a VM with a PCI-attached TPU board with four dual-core TPU chips on it.</p> </div> <p><a href="https://codelabs.developers.google.com/codelabs/keras-flowers-convnets/#2">https://codelabs.developers.google.com/codelabs/keras-flowers-convnets/#2</a></p> <div data-bbox="970 760 1696 1182">  <p>TPU v2 - 4 chips, 2 cores per chip</p> </div> <p><a href="https://cloud.google.com/tpu/docs/system-architecture">https://cloud.google.com/tpu/docs/system-architecture</a></p>

'273 PATENT	INFRINGEMENT EVIDENCE
<p>53. A <b>device</b>:  comprising at least one first low precision high-dynamic range (LPHDR) execution unit adapted to execute a first operation on a first input signal representing a first numerical value to produce a first output signal representing a second numerical value,  wherein the dynamic range of the possible valid inputs to the first operation is at least as wide as from 1/1,000,000 through 1,000,000 and for at least X=5% of the possible valid inputs to the first operation, the statistical mean, over repeated execution of the first operation on each specific input from the at least X% of the possible valid inputs to the first operation, of the numerical values represented by the first output signal of the LPHDR unit executing the first operation on that input differs by at least Y=0.05% from the result of an exact mathematical calculation of the first operation on the numerical values of that same input;  wherein the number of LPHDR execution units in the device exceeds by at least one hundred the non-negative integer number of execution units in the device adapted to execute at least the operation of multiplication on floating point numbers that are at least 32 bits wide.</p>	<p>As demonstrated below, the Accused Products include multiple components that, separately and independently, meet all the requirements of the claimed "device." For example, a "TPU Chip" satisfies these requirements:</p> <div data-bbox="919 272 1738 743"> <p>When you request one "Cloud TPU v2" on Google Cloud Platform, you get a virtual machine (VM) which has a PCI-attached TPU board. The TPU board has four dual-core TPU chips. Each TPU core features a VPU (Vector Processing Unit) and a 128x128 MXU (Matrix multiply Unit). This "Cloud TPU" is then usually connected through the network to the VM that requested it. So the full picture looks like this:</p>  <p>Illustration: your VM with a network-attached "Cloud TPU" accelerator. "The Cloud TPU" itself is made of a VM with a PCI-attached TPU board with four dual-core TPU chips on it.</p> </div> <p><a href="https://codelabs.developers.google.com/codelabs/keras-flowers-convnets/#2">https://codelabs.developers.google.com/codelabs/keras-flowers-convnets/#2</a></p> <div data-bbox="982 813 1682 1230">  <p>TPU v2 - 4 cores, 2 cores per chip  TPU v3 - 4 cores, 2 cores per chip</p> </div> <p><a href="https://cloud.google.com/tpu/docs/system-architecture">https://cloud.google.com/tpu/docs/system-architecture</a></p> <p>See also generally Norrie et al., "Google's Training Chips Revealed: TPuv2 and TPuv3" (Presented at HotChips Conference, Aug. 2020)</p>



## '273 PATENT

53. A **device**:

comprising at least one first low precision high-dynamic range (LPHDR) execution unit adapted to execute a first operation on a first input signal representing a first numerical value to produce a first output signal representing a second numerical value,

wherein the dynamic range of the possible valid inputs to the first operation is at least as wide as from 1/1,000,000 through 1,000,000 and for at least  $X=5\%$  of the possible valid inputs to the first operation, the statistical mean, over repeated execution of the first operation on each specific input from the at least  $X\%$  of the possible valid inputs to the first operation, of the numerical values represented by the first output signal of the LPHDR unit executing the first operation on that input differs by at least  $Y=0.05\%$  from the result of an exact mathematical calculation of the first operation on the numerical values of that same input;

wherein the number of LPHDR execution units in the device exceeds by at least one hundred the non-negative integer number of execution units in the device adapted to execute at least the operation of multiplication on floating point numbers that are at least 32 bits wide.

## INFRINGEMENT EVIDENCE

As demonstrated below, the Accused Products include multiple components that, separately and independently, meet all the requirements of the claimed "device." For example, a "TPU Core" satisfies these requirements:

When you request one "Cloud TPU v2" on Google Cloud Platform, you get a virtual machine (VM) which has a PCI-attached TPU board. The TPU board has four dual-core TPU chips. Each TPU core features a VPU (Vector Processing Unit) and a 128x128 MXU (Matrix multiply Unit). This "Cloud TPU" is then usually connected through the network to the VM that requested it. So the full picture looks like this:

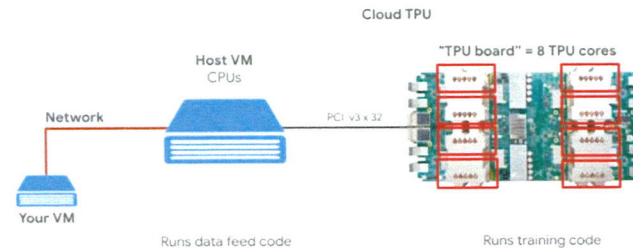
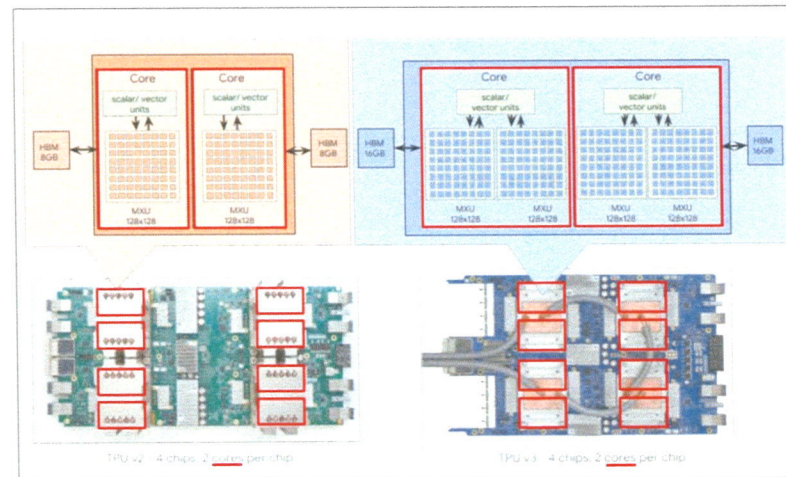


Illustration: your VM with a network-attached "Cloud TPU" accelerator. "The Cloud TPU" itself is made of a VM with a PCI-attached TPU board with four dual-core TPU chips on it.

<https://codelabs.developers.google.com/codelabs/keras-flowers-convnets/#2>



<https://cloud.google.com/tpu/docs/system-architecture>

## '273 PATENT

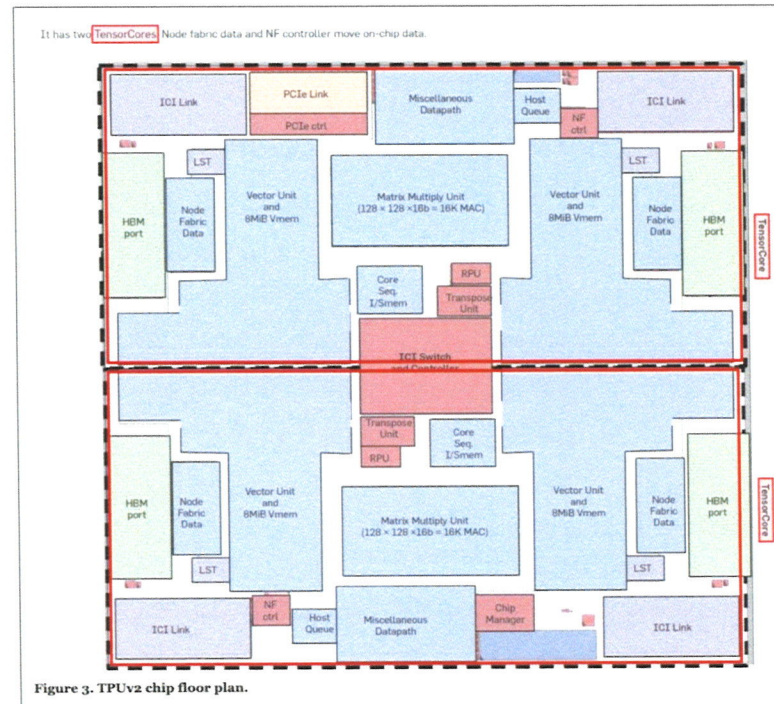
53. A **device**:

comprising at least one first low precision high-dynamic range (LPHDR) execution unit adapted to execute a first operation on a first input signal representing a first numerical value to produce a first output signal representing a second numerical value,

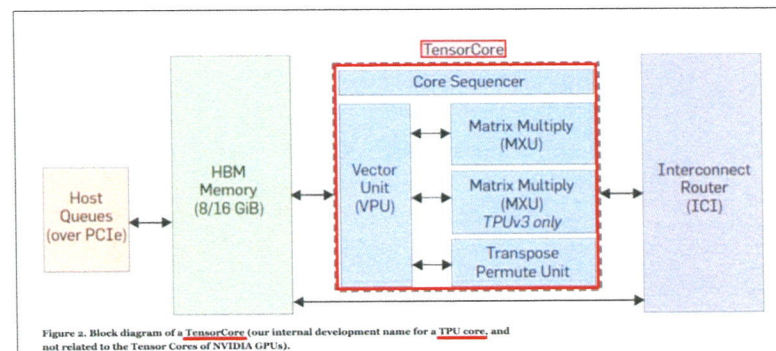
wherein the dynamic range of the possible valid inputs to the first operation is at least as wide as from 1/1,000,000 through 1,000,000 and for at least  $X=5\%$  of the possible valid inputs to the first operation, the statistical mean, over repeated execution of the first operation on each specific input from the at least  $X\%$  of the possible valid inputs to the first operation, of the numerical values represented by the first output signal of the LPHDR unit executing the first operation on that input differs by at least  $Y=0.05\%$  from the result of an exact mathematical calculation of the first operation on the numerical values of that same input;

wherein the number of LPHDR execution units in the device exceeds by at least one hundred the non-negative integer number of execution units in the device adapted to execute at least the operation of multiplication on floating point numbers that are at least 32 bits wide.

## INFRINGEMENT EVIDENCE



<https://cacm.acm.org/magazines/2020/7/245702-a-domain-specific-supercomputer-for-training-deep-neural-networks>



*Id.*



## '273 PATENT

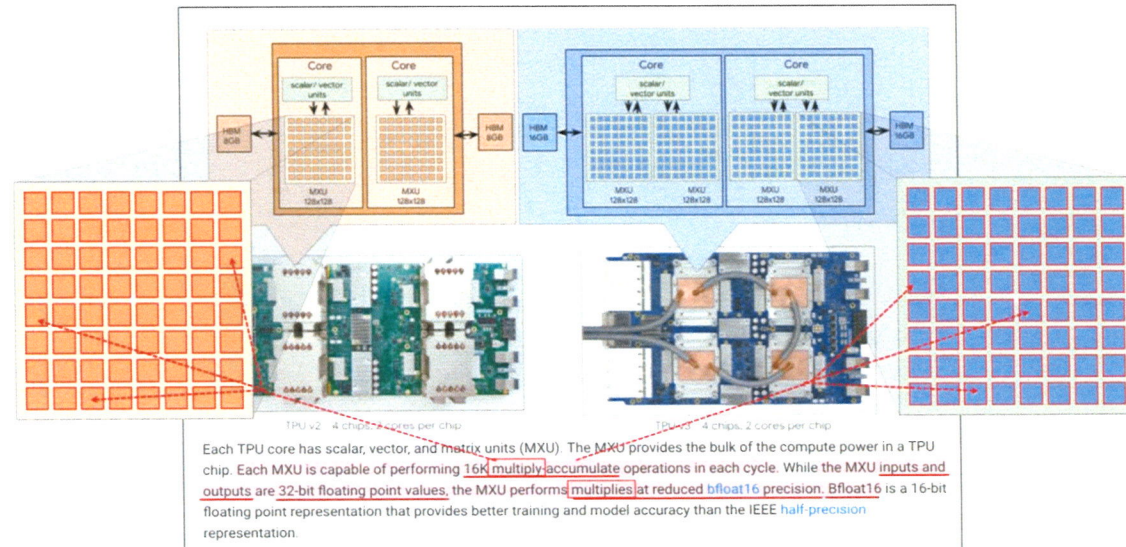
53. A device:

comprising at least one first low precision high-dynamic range (LPHDR) execution unit adapted to execute a first operation on a first input signal representing a first numerical value to produce a first output signal representing a second numerical value,

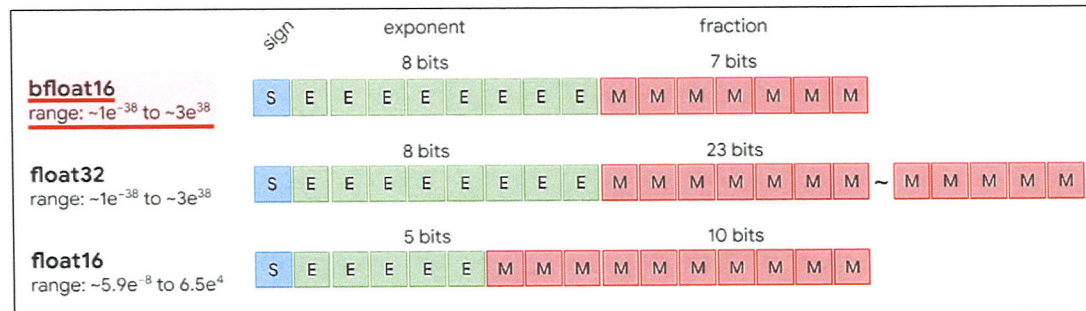
wherein the dynamic range of the possible valid inputs to the first operation is at least as wide as from 1/1,000,000 through 1,000,000 and for at least  $X=5\%$  of the possible valid inputs to the first operation, the statistical mean, over repeated execution of the first operation on each specific input from the at least  $X\%$  of the possible valid inputs to the first operation, of the numerical values represented by the first output signal of the LPHDR unit executing the first operation on that input differs by at least  $Y=0.05\%$  from the result of an exact mathematical calculation of the first operation on the numerical values of that same input;

wherein the number of LPHDR execution units in the device exceeds by at least one hundred the non-negative integer number of execution units in the device adapted to execute at least the operation of multiplication on floating point numbers that are at least 32 bits wide.

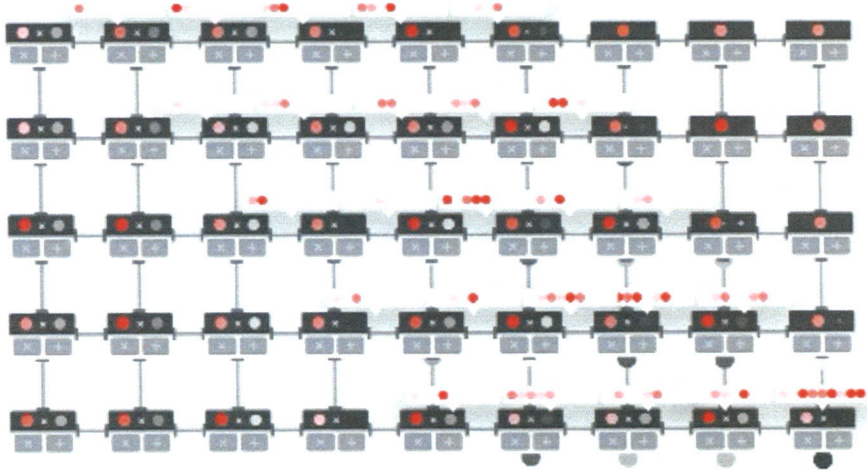
## INFRINGEMENT EVIDENCE



<https://cloud.google.com/tpu/docs/system-architecture>

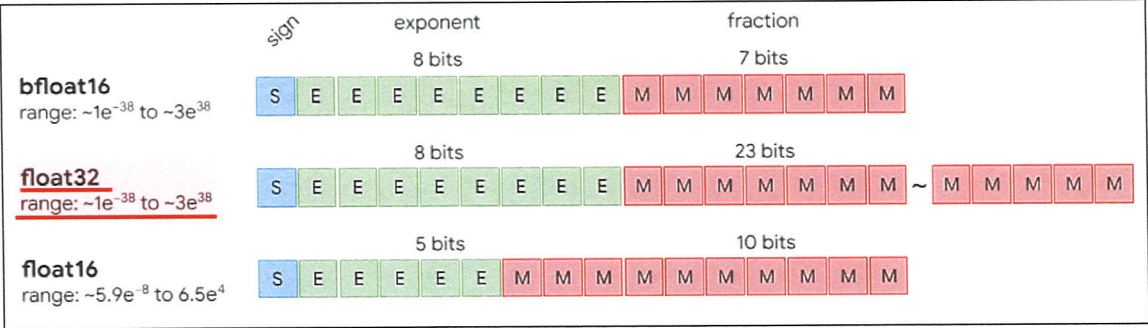


<https://cloud.google.com/tpu/docs/bfloat16>

'273 PATENT	INFRINGEMENT EVIDENCE
<p>53. A device:</p> <p><b>comprising at least one first low precision high-dynamic range (LPHDR) execution unit adapted to execute a first operation on a first input signal representing a first numerical value to produce a first output signal representing a second numerical value,</b></p> <p>wherein the dynamic range of the possible valid inputs to the first operation is at least as wide as from 1/1,000,000 through 1,000,000 and for at least X=5% of the possible valid inputs to the first operation, the statistical mean, over repeated execution of the first operation on each specific input from the at least X% of the possible valid inputs to the first operation, of the numerical values represented by the first output signal of the LPHDR unit executing the first operation on that input differs by at least Y=0.05% from the result of an exact mathematical calculation of the first operation on the numerical values of that same input;</p> <p>wherein the number of LPHDR execution units in the device exceeds by at least one hundred the non-negative integer number of execution units in the device adapted to execute at least the operation of multiplication on floating point numbers that are at least 32 bits wide.</p>	<p><b>Systolic array</b></p> <p>The MXU implements matrix multiplications in hardware using a so-called "systolic array" architecture in which <u>data elements flow through an array of hardware computation units.</u> (In medicine, "systolic" refers to heart contractions and blood flow, here to the flow of data.)</p> <p>The basic element of a matrix multiplication is a dot product between a line from one matrix and a column from the other matrix (see illustration at the top of this section). For a matrix multiplication <math>Y=X*W</math>, one element of the result would be:</p> $Y[2,0] = X[2,0]*W[0,0] + X[2,1]*W[1,0] + X[2,2]*W[2,0] + \dots + X[2,n]*W[n,0]$  <p><i>Illustration: the MXU systolic array. The compute elements are multiply-accumulators. <u>The values of one matrix are loaded into the array (red dots).</u> <u>Values of the other matrix flow through the array (grey dots).</u> Vertical lines propagate the values up. Horizontal lines propagate partial sums. It is left as an exercise to the user to verify that as the data flows through the array, you get the result of the matrix multiplication coming out of the right side.</i></p> <p><a href="https://codelabs.developers.google.com/codelabs/keras-flowers-convnets/#2">https://codelabs.developers.google.com/codelabs/keras-flowers-convnets/#2</a></p>



'273 PATENT	INFRINGEMENT EVIDENCE
<p>53. A device:  comprising at least one first low precision high-dynamic range (LPHDR) execution unit adapted to execute a first operation on a first input signal representing a first numerical value to produce a first output signal representing a second numerical value,</p> <p><b>wherein the dynamic range of the possible valid inputs to the first operation is at least as wide as from 1/1,000,000 through 1,000,000</b> and for at least X=5% of the possible valid inputs to the first operation, the statistical mean, over repeated execution of the first operation on each specific input from the at least X% of the possible valid inputs to the first operation, of the numerical values represented by the first output signal of the LPHDR unit executing the first operation on that input differs by at least Y=0.05% from the result of an exact mathematical calculation of the first operation on the numerical values of that same input;</p> <p>wherein the number of LPHDR execution units in the device exceeds by at least one hundred the non-negative integer number of execution units in the device adapted to execute at least the operation of multiplication on floating point numbers that are at least 32 bits wide.</p>	<ul style="list-style-type: none"> <li>• “Each TPU core has scalar, vector, and matrix units (MXU). The MXU provides the bulk of the compute power in a TPU chip. Each MXU is capable of performing 16K multiply-accumulate operations in each cycle. While <b>the MXU inputs and outputs are 32-bit floating point values</b>, the MXU performs multiplies at reduced bfloat16 precision. Bfloat16 is a 16-bit floating point representation that provides better training and model accuracy than the IEEE half-precision representation.”  <a href="https://cloud.google.com/tpu/docs/system-architecture">https://cloud.google.com/tpu/docs/system-architecture</a></li> <li>• “The following figure shows three floating-point[] formats <ul style="list-style-type: none"> <li>• <b>fp32 - IEEE single-precision floating-point</b></li> <li>• fp16 - IEEE half-precision floating point</li> <li>• bfloat16 - 16-bit <i>brain floating point</i>”</li> </ul> <a href="https://cloud.google.com/tpu/docs/bfloat16">https://cloud.google.com/tpu/docs/bfloat16</a></li> </ul> <div data-bbox="787 544 1869 852"> <p>The diagram illustrates the bit layouts for three floating-point formats:</p> <ul style="list-style-type: none"> <li><b>bfloat16</b>: range: <math>\sim 1e^{-38}</math> to <math>\sim 3e^{38}</math>. It consists of a sign bit (S), an 8-bit exponent (E), and a 7-bit fraction (M).</li> <li><b>float32</b>: range: <math>\sim 1e^{-38}</math> to <math>\sim 3e^{38}</math>. It consists of a sign bit (S), an 8-bit exponent (E), and a 23-bit fraction (M).</li> <li><b>float16</b>: range: <math>\sim 5.9e^{-8}</math> to <math>6.5e^4</math>. It consists of a sign bit (S), a 5-bit exponent (E), and a 10-bit fraction (M).</li> </ul> </div> <p><i>Id.</i></p>

'273 PATENT	INFRINGEMENT EVIDENCE
<p>53. A device:          comprising at least one first low precision high-dynamic range (LPHDR) execution unit adapted to execute a first operation on a first input signal representing a first numerical value to produce a first output signal representing a second numerical value,          wherein the dynamic range of the possible valid inputs to the first operation is at least as wide as from 1/1,000,000 through 1,000,000 and <b>for at least X=5% of the possible valid inputs to the first operation, the statistical mean, over repeated execution of the first operation on each specific input from the at least X% of the possible valid inputs to the first operation, of the numerical values represented by the first output signal of the LPHDR unit executing the first operation on that input differs by at least Y=0.05% from the result of an exact mathematical calculation of the first operation on the numerical values of that same input;</b>          wherein the number of LPHDR execution units in the device exceeds by at least one hundred the non-negative integer number of execution units in the device adapted to execute at least the operation of multiplication on floating point numbers that are at least 32 bits wide.</p>	<ul style="list-style-type: none"> <li>“Each TPU core has scalar, vector, and matrix units (MXU). The MXU provides the bulk of the compute power in a TPU chip. Each MXU is capable of performing 16K multiply-accumulate operations in each cycle. While <b>the MXU inputs and outputs are 32-bit floating point values</b>, the MXU performs multiplies at reduced bfloat16 precision. Bfloat16 is a 16-bit floating point representation that provides better training and model accuracy than the IEEE half-precision representation.”  <a href="https://cloud.google.com/tpu/docs/system-architecture">https://cloud.google.com/tpu/docs/system-architecture</a></li> <li>“The following figure shows three floating-point[] formats             <ul style="list-style-type: none"> <li><b>fp32 - IEEE single-precision floating-point</b></li> <li>fp16 - IEEE half-precision floating point</li> <li>bfloat16 - 16-bit <i>brain floating point</i>”</li> </ul> <a href="https://cloud.google.com/tpu/docs/bfloat16">https://cloud.google.com/tpu/docs/bfloat16</a> </li> </ul>  <p><i>Id.</i></p> <ul style="list-style-type: none"> <li>“Because general-purpose processors such as CPUs and GPUs must provide good performance across a wide range of applications, they have evolved myriad sophisticated, performance-oriented mechanisms. As a side effect, the behavior of those processors can be difficult to predict, which makes it hard to guarantee a certain latency limit on neural network inference. In contrast, TPU design is strictly <b>minimal and deterministic</b> as it has to run <b>only one task at a time: neural network prediction</b>. You can see its simplicity in the floor plan of the TPU die.”  <a href="https://cloud.google.com/blog/products/gcp/an-in-depth-look-at-googles-first-tensor-processing-unit-tpu">https://cloud.google.com/blog/products/gcp/an-in-depth-look-at-googles-first-tensor-processing-unit-tpu</a> (emphasis in orig.)</li> <li>“In mathematics, computer science and physics, a <b>deterministic</b> system is a system in which <b>no randomness</b> is involved in the development of future states of the system. A <b>deterministic model will thus always produce the same output from a given starting condition or initial state.</b>”  <a href="https://en.wikipedia.org/wiki/Deterministic_system">https://en.wikipedia.org/wiki/Deterministic_system</a></li> <li>For each of the possible valid inputs to the multiplication operation performed by the multipliers within the MXU, Singular has computed the result and compared it to the result of an exact mathematical calculation performed on the same inputs. The results of this test showed that for <b>more than 10%</b> of the possible valid inputs, the numerical value represented by the output signal of each MXU multiplier differs by <b>more than 0.2%</b> from the result of an exact mathematical calculation performed on the same inputs.</li> </ul>



## '273 PATENT

53. A **device**:

comprising at least one first low precision high-dynamic range (LPHDR) execution unit adapted to execute a first operation on a first input signal representing a first numerical value to produce a first output signal representing a second numerical value,

wherein the dynamic range of the possible valid inputs to the first operation is at least as wide as from 1/1,000,000 through 1,000,000 and for at least X=5% of the possible valid inputs to the first operation, the statistical mean, over repeated execution of the first operation on each specific input from the at least X% of the possible valid inputs to the first operation, of the numerical values represented by the first output signal of the LPHDR unit executing the first operation on that input differs by at least Y=0.05% from the result of an exact mathematical calculation of the first operation on the numerical values of that same input;

**wherein the number of LPHDR execution units in the device exceeds by at least one hundred the non-negative integer number of execution units in the device adapted to execute at least the operation of multiplication on floating point numbers that are at least 32 bits wide.**

## INFRINGEMENT EVIDENCE

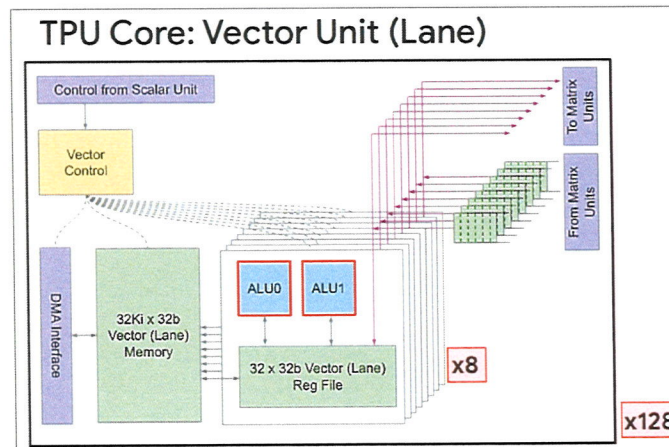
The Accused Products independently meet this claim limitation for each "device" identified above:

We cannot reveal technology details of our chip partner. Although it is a larger, older technology, the TPUv2 die size is less than 3/4s of the GPU. TPUv3 is 6% larger in that same technology. TDP stands for Thermal Design Power. The Volta has 80 symmetric multiprocessors.

Feature	TPUv1	TPUv2	TPUv3	Volta
Peak TeraFLOPS/Chip	92 (8b int)	46 (16b) 3 (32b)	123 (16b) 4 (32b)	125 (16b) 16 (32b)
Network links x Gbits/s/Chip	—	4 x 496	4 x 656	6 x 200
Max chips/supercomputer	—	256	1024	Varies
Peak PetaFLOPS/supercomputer	—	11.8	126	Varies
Bisection Terabits/supercomputer	—	15.9	42.0	Varies
Clock Rate (MHz)	700	700	940	1530
TDP (Watts)/Chip	75	280	450	450
TDP (Kwatts)/supercomputer	—	124	594	Varies
Die Size (mm <sup>2</sup> )	<331	<611	<648	815
Chip Technology	28nm	>12nm	>12nm	12nm
Memory size (on/off-chip)	28MB/8GB	32MB/16GB	32MB/32GB	36MB/32GB
Memory GB/s/Chip	34	700	900	900
MXUs/Core	1	128x128	128x128	8
MXU Size	256x256			4x4
Cores/Chip	1	2	2	80
Chips/CPU Host	4	4	8	8 or 16

**Table 3. Key processor features.**

<https://cacm.acm.org/magazines/2020/7/245702-a-domain-specific-supercomputer-for-training-deep-neural-networks/fulltext>



Norrie et al., "Google's Training Chips Revealed: TPUv2 and TPUv3"  
(Presented at HotChips Conference, Aug. 2020)